

## 5.4 The chi-squared test

You may be interested in finding out whether or not certain sets of data are independent. Suppose you collect data on the favorite color of T-shirt for men and women. You may want to find out whether color and gender are independent or not. One way to do this is to perform a **chi-squared test** ( $\chi^2$ ) for independence.

To perform a chi-squared test ( $\chi^2$ ) there are four main steps.

**Step 1:** Write the **null** ( $H_0$ ) and **alternative** ( $H_1$ ) hypotheses.

$H_0$  states that the data sets are independent.

$H_1$  states that the data sets are not independent.

For example, the hypotheses for color of T-shirt and gender could be:

$H_0$ : Color of T-shirt is independent of gender.

$H_1$ : Color of T-shirt is not independent of gender.

**Step 2:** Calculate the chi-squared test statistic.

Firstly, you may need to put the data into a **contingency table**, which shows the frequencies of two variables. The elements in the table are the **observed** data. The elements should be frequencies (not percentages).

For the example above, the contingency table could be:

	Black	White	Red	Blue	Totals
Male	48	12	33	57	150
Female	35	46	42	27	150
Totals	83	58	75	84	300

If you are given the contingency table, you may need to extend it to include an extra row and column for the 'Totals'.

From the observed data, you can calculate the **expected frequencies**. Since you are testing for independence, you can use the formula for the probability of independent events to calculate the expected values. So:

The expected number of men who like black T-shirts is

$$\frac{150}{300} \times \frac{83}{300} \times 300 = 41.5.$$

The expected number of men who like white T-shirts is

$$\frac{150}{300} \times \frac{58}{300} \times 300 = 29 \text{ and so on.}$$

The expected table of values would then look like this:

	Black	White	Red	Blue	Totals
Male	41.5	29	37.5	42	150
Female	41.5	29	37.5	42	150
Totals	83	58	75	84	300

When two variables are independent, one does not affect the other. Here, you are finding out whether a person's gender influences their colour choice. You will learn more about mathematical independence in Chapter 8.

The main entries in this table form a  $2 \times 4$  **matrix** (array of numbers) - do not include the row and column for the totals.

In examinations, the **largest contingency table** will be a  $4 \times 4$ .

### Note:

- The expected values can **never** be less than 1.
- The expected values must be 5 or higher.
- If there are entries between 1 and 5, you can combine table rows or columns.

For calculations by hand, you need the expected frequencies to find the  $\chi^2$  value.

→ To calculate the  $\chi^2$  value use the formula  $\chi^2_{\text{calc}} = \sum \frac{(f_o - f_e)^2}{f_e}$ , where  $f_o$  are the observed frequencies and  $f_e$  are the expected frequencies.

For our example,

$$\begin{aligned}\chi^2_{\text{calc}} &= \frac{(48-41.5)^2}{41.5} + \frac{(12-29)^2}{29} + \frac{(33-37.5)^2}{37.5} + \frac{(57-42)^2}{42} + \frac{(35-41.5)^2}{41.5} \\ &\quad + \frac{(46-29)^2}{29} + \frac{(42-37.5)^2}{37.5} + \frac{(27-42)^2}{42} \\ &= 33.8\end{aligned}$$

Using your GDC to find the  $\chi^2$  value, enter the contingency table as a matrix (array) and then use the matrix with the  $\chi^2$  2-way test.

State Apps			
48	12	33	57
35	46	42	27

→mat

χ² 2way mat: stat results	
"Title"	"χ² 2-way Test"
"χ²"	33.7615
"PVal"	2.22473E-7
"df"	3.
"ExpMatrix"	"[...]"
"CompMatrix"	"[...]"

From the screenshot, you can see that  $\chi^2_{\text{calc}} = 33.8$  (to 3 sf). This confirms our earlier hand calculation.

**Step 3:** Calculate the critical value.

First note the **level of significance**. This is given in examination questions but you have to decide which level to use in your project. The most common levels are 1%, 5% and 10%.

Now you need to calculate the number of **degrees of freedom**.

→ To find the degrees of freedom for the chi-squared test for independence, use this formula based on the contingency table.

$$\text{Degrees of freedom} = (\text{number of rows} - 1) \times (\text{number of columns} - 1)$$

If the number of degrees of freedom is 1, you will be expected to use **Yates' continuity correction** to work out the chi-squared value. (In examinations the degrees of freedom will always be greater than 1.)

So, in our ongoing example, the number of degrees of freedom is  $(2 - 1) \times (4 - 1) = 3$

Your GDC calculates the expected values for you but you must know how to find them by hand in case you are asked to show one or two calculations in an exam question. To see the matrix for the expected values, type 'stat.' and then select 'expmatrix' from the menu that pops up.

GDC help on CD: Alternative demonstrations for the TI-84 Plus and Casio FX-9860GII GDCs are on the CD.



The level of significance and degrees of freedom can be used to find the critical value. However, in examinations, the **critical value** will always be given.

For our example, at the 1% level, the critical value is 11.345. At the 5% level, the critical value is 7.815. At the 10% level, the critical value is 6.251.

**Step 4:** Compare  $\chi^2_{\text{calc}}$  against the critical value.

→ If  $\chi^2_{\text{calc}}$  is **less than** the critical value then **do not reject** the null hypothesis.

If  $\chi^2_{\text{calc}}$  is **more than** the critical value then **reject** the null hypothesis.

In our example, at the 5% level,  $33.8 > 7.815$ . Therefore, we reject the null hypothesis that T-shirt color is independent of gender.

Using a GDC, you can compare the  $p$ -value against the significance level.

→ If the  $p$ -value is **less** than the significance level then **reject** the null hypothesis.

If the  $p$ -value is **more** than the significance level then **do not reject** the null hypothesis.

The  $p$ -value is the probability value. It is the probability of evidence against the null hypothesis.

Use the significance level as a decimal, so 1% = 0.01, 5% = 0.05 and 10% = 0.1.

So, for our example,  $p\text{-value} = 0.000\,000\,2$  (see the GDC screenshot on page 234).

$0.000\,000\,2 < 0.05$ , so we reject the null hypothesis.

→ **To perform a  $\chi^2$  test:**

- 1 Write the null ( $H_0$ ) and alternative ( $H_1$ ) hypotheses.
- 2 Calculate  $\chi^2_{\text{calc}}$ :
  - a using your GDC (examinations)
  - b using the  $\chi^2_{\text{calc}}$  formula (project work)
- 3 Determine:
  - a the  $p$ -value by using your GDC
  - b the critical value (given in examinations)
- 4 Compare:
  - a the  $p$ -value against the significance level
  - b  $\chi^2_{\text{calc}}$  against the critical value

### Example 13

One hundred people were interviewed outside a chocolate shop to find out which flavor of chocolate cream they preferred. The results are given in the table, classified by gender.

	Strawberry	Coffee	Orange	Vanilla	Totals
Male	23	18	8	8	57
Female	15	6	12	10	43
Totals	38	24	20	18	100

Perform a  $\chi^2$  test, at the 5% significance level, to determine whether the flavor of chocolate cream is independent of gender.

- State the null hypothesis and the alternative hypothesis.
- Show that the expected frequency for female and strawberry flavor is approximately 16.3.
- Write down the number of degrees of freedom.
- Write down the  $\chi^2_{\text{calc}}$  value for this data.

The critical value is 7.815.

- Using the critical value or the  $p$ -value, comment on your result.

#### Answers

- $H_0$ : Flavor of chocolate cream is independent of gender.  
 $H_1$ : Flavor of chocolate cream is not independent of gender.

$$\text{b } \frac{43}{100} \times \frac{38}{100} \times 100 = 16.34$$

So, the expected frequency for female and strawberry flavor is approximately 16.3.

$$\text{c } \text{Degrees of freedom} = (2 - 1)(4 - 1) = 3$$

$$\text{d } \chi^2_{\text{calc}} = 6.88$$

- $6.88 < 7.815$ ; therefore, we do not reject the null hypothesis. There is enough evidence to conclude that flavor of chocolate cream is independent of gender.

Write  $H_0$  using 'independent of'.

Write  $H_1$  using 'not independent of'.

From the contingency table:

Total for 'female' row = 43

Total for 'strawberry' column = 38

Total surveyed = 100

Degrees of freedom = (number of rows - 1) (number of columns - 1)

Here, there are 2 rows and 4 columns in the observed matrix of the contingency table.

Using your GDC:

Enter the contingency table as a matrix. Use the matrix with  $\chi^2$  2-way test. Read off  $\chi^2$  value.

The  $p$ -value = 0.0758.

Using the given critical value, check:

$\chi^2_{\text{calc}} < \text{critical value} \rightarrow \text{do not reject, or}$

$\chi^2_{\text{calc}} > \text{critical value} \rightarrow \text{reject.}$

Or, using the  $p$ -value, check:

$p\text{-value} < \text{significance level} \rightarrow \text{reject, or}$

$p\text{-value} > \text{significance level} \rightarrow \text{do not reject.}$

Significance level = 5% = 0.05. So,  $0.0758 > 0.05$  and we do not reject the null hypothesis.